# Advanced Data Structures and Algorithm Analysis

丁尧相

浙江大学

Fall and Winter 2025
Lecture 8

# Inverted File Index

- Inverted File Index

- Take-home messages

# Inverted File Index

- **Inverted File Index**

- Take-home messages

# AVL Trees

# AVL Trees

## Splay Trees

# AVL Trees

## Splay Trees

### Red-Black Trees

**AVL Trees**

**Splay Trees**

**Red-Black Trees**

**B+ Trees**

# AVL Trees

## Splay Trees

### Red-Black Trees

### B+ Trees

### 1 | Google

1 - eBizMBA Rank | **1,800,000,000** - Estimated Unique Monthly Visitors | 1 - Quantcast Rank | 1 - Alexa Rank | 1 - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA

### 2 | Bing

33 - eBizMBA Rank | **500,000,000** - Estimated Unique Monthly Visitors | 8 - Quantcast Rank | 40 - Alexa Rank | 43 - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA

### 3 | Yahoo! Search

43 - eBizMBA Rank | **490,000,000** - Estimated Unique Monthly Visitors | 8 - Quantcast Rank | *56* - Alexa Rank | *67* - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA

### 4 | Baidu

54 - eBizMBA Rank | **480,000,000** - Estimated Unique Monthly Visitors | *150* - Quantcast Rank | 4 - Alexa Rank | 9 - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA
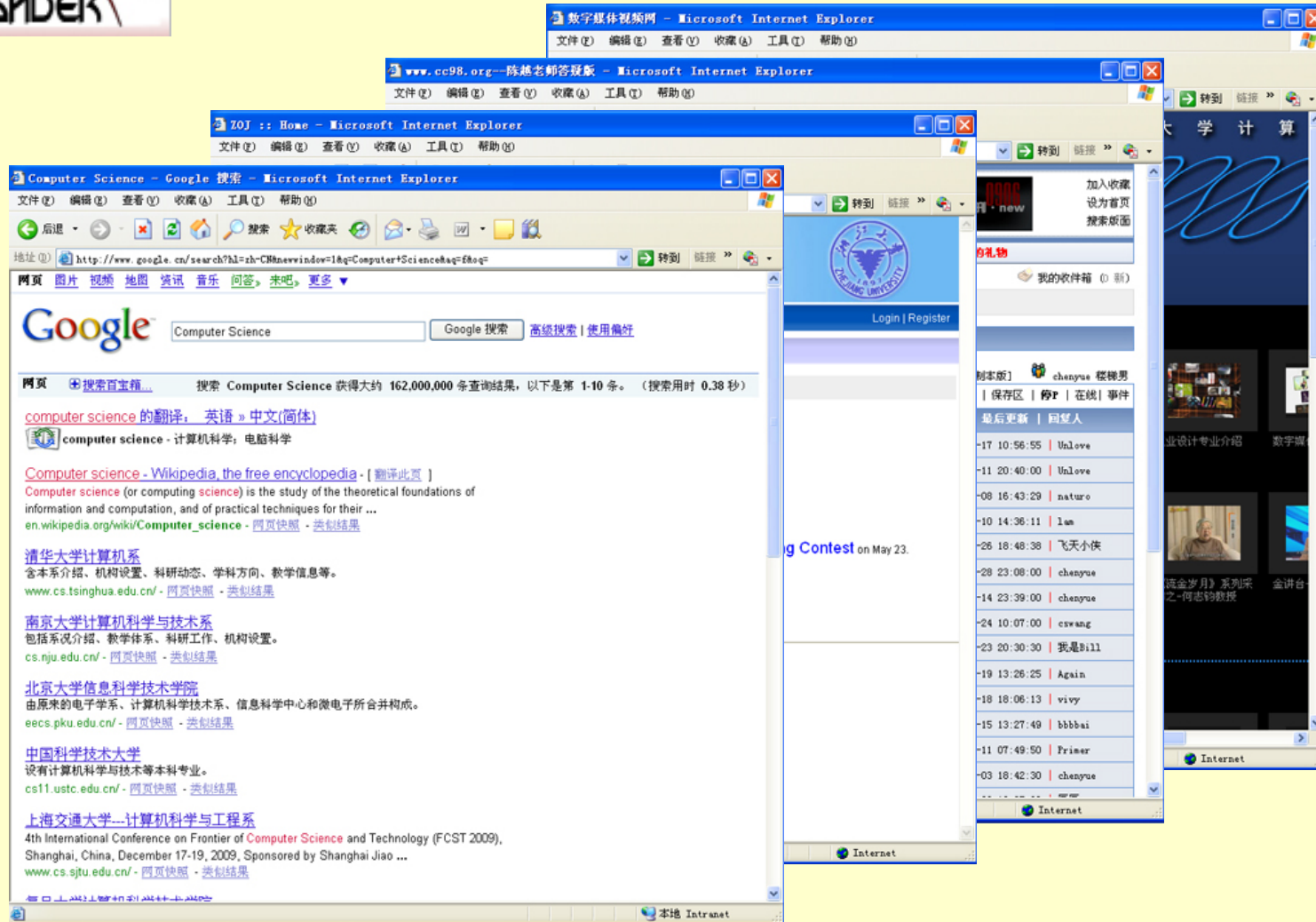
### 5 | Ask

205 - eBizMBA Rank | **300,000,000** - Estimated Unique Monthly Visitors | 329 - Quantcast Rank | 110 - Alexa Rank | 177 - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA

### 6 | Aol Search

273 - eBizMBA Rank | **200,000,000** - Estimated Unique Monthly Visitors | *350* - Quantcast Rank | 276 - Alexa Rank | *194* - SimilarWeb Rank | *Last Updated:* February 1, 2020. The Best Search Engines | eBizMBA

How can I find in which retrieved web pages that include "*Computer Science*"?

6

# Information Retrieval

# Information Retrieval

- Information Retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).

# Information Retrieval

- Information Retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).

  - These days we frequently think first of web search, but there are many other cases:
    - E-mail search
    - Searching your laptop
    - Corporate knowledge bases
    - Legal information retrieval

# Basic assumptions of Information Retrieval
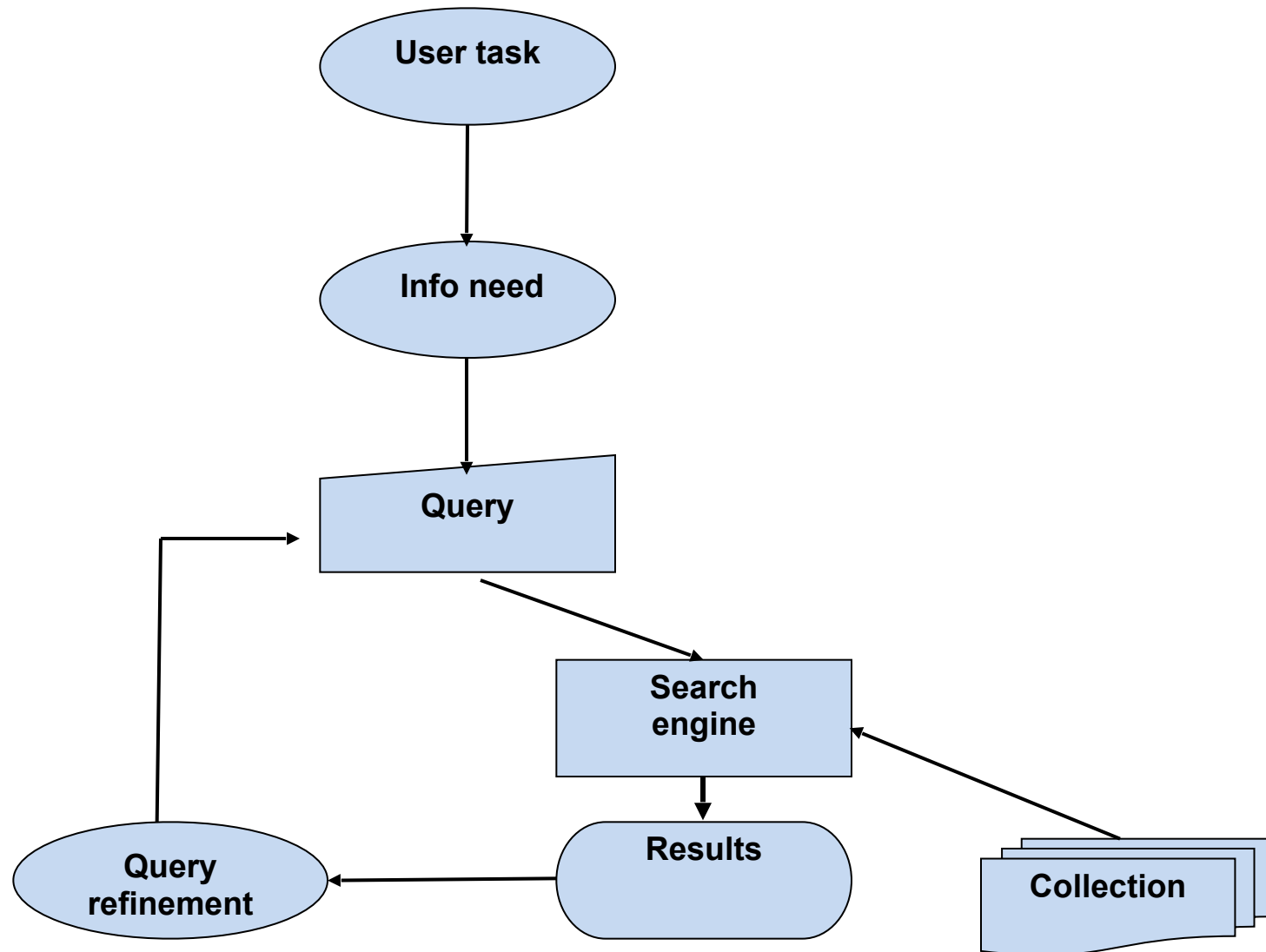
# Basic assumptions of Information Retrieval

- Collection: A set of documents
  - Assume it is a static collection for the moment
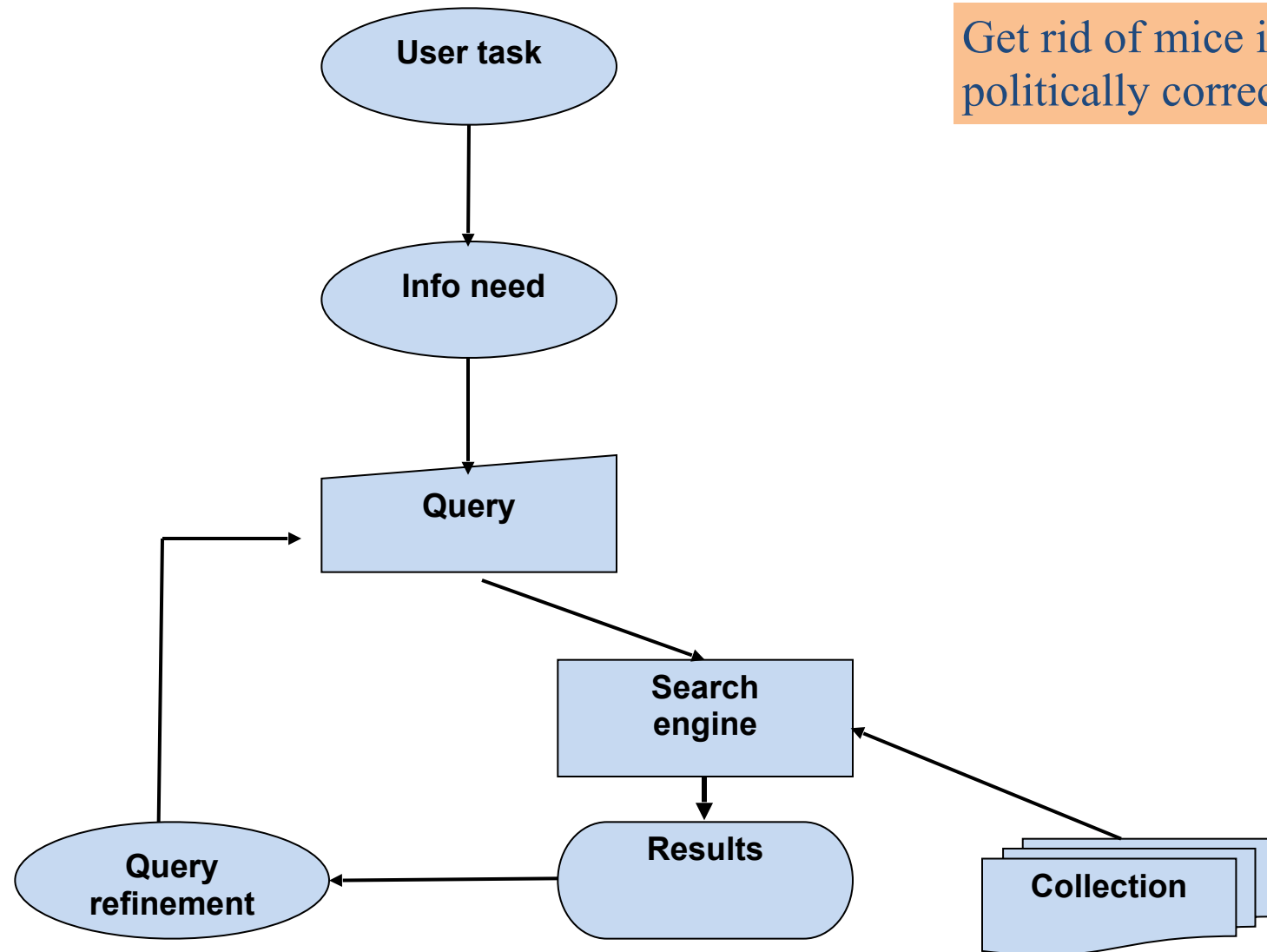
# Basic assumptions of Information Retrieval

- Collection: A set of documents
  - Assume it is a static collection for the moment

- Goal: Retrieve documents with information that is relevant to the user's information need and helps the user complete a task

This is different from searching for a fixed key in a data structure or database

# The classic search model

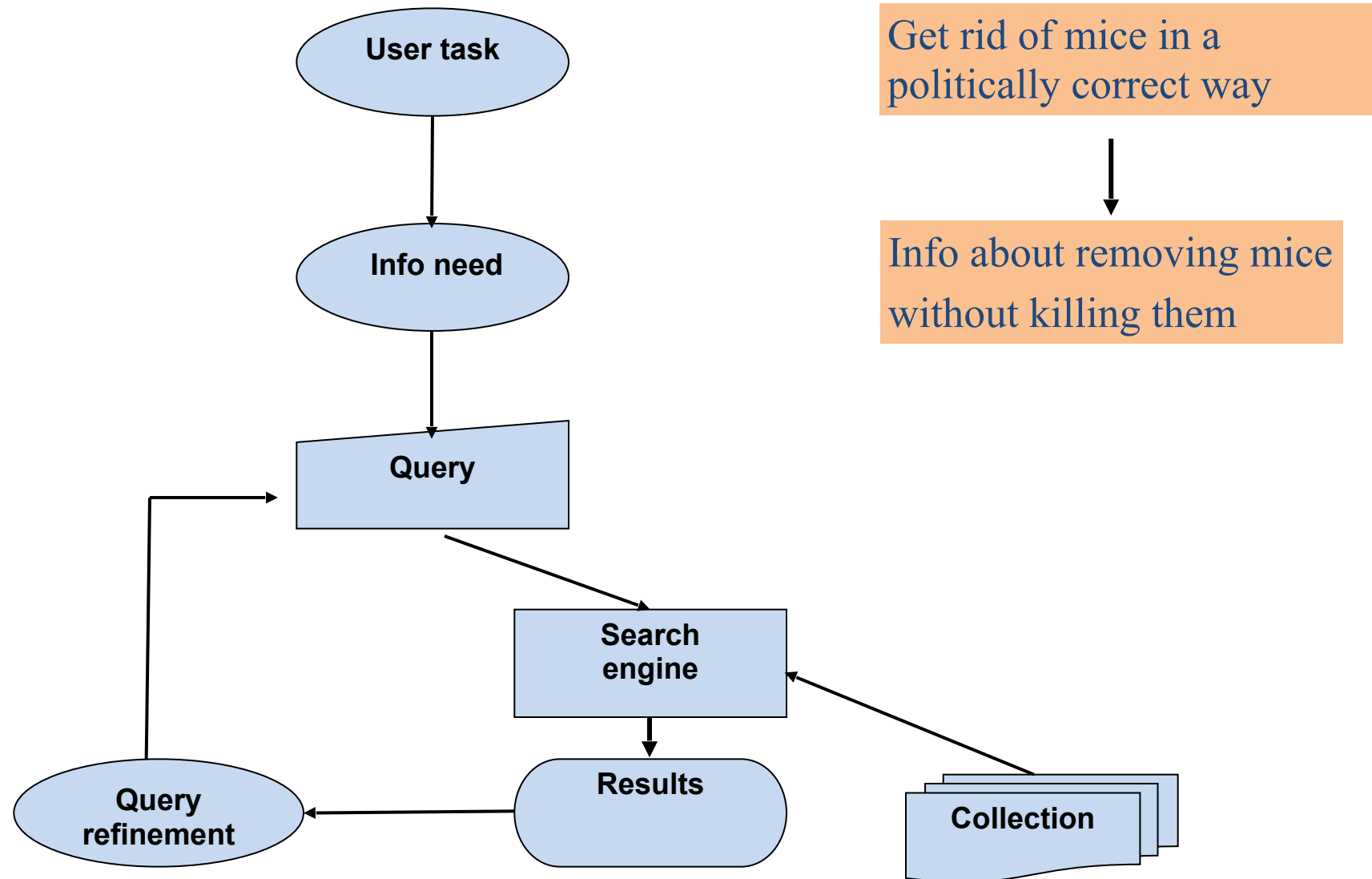# The classic search model



User task

Info need

Query

Search engine

Results
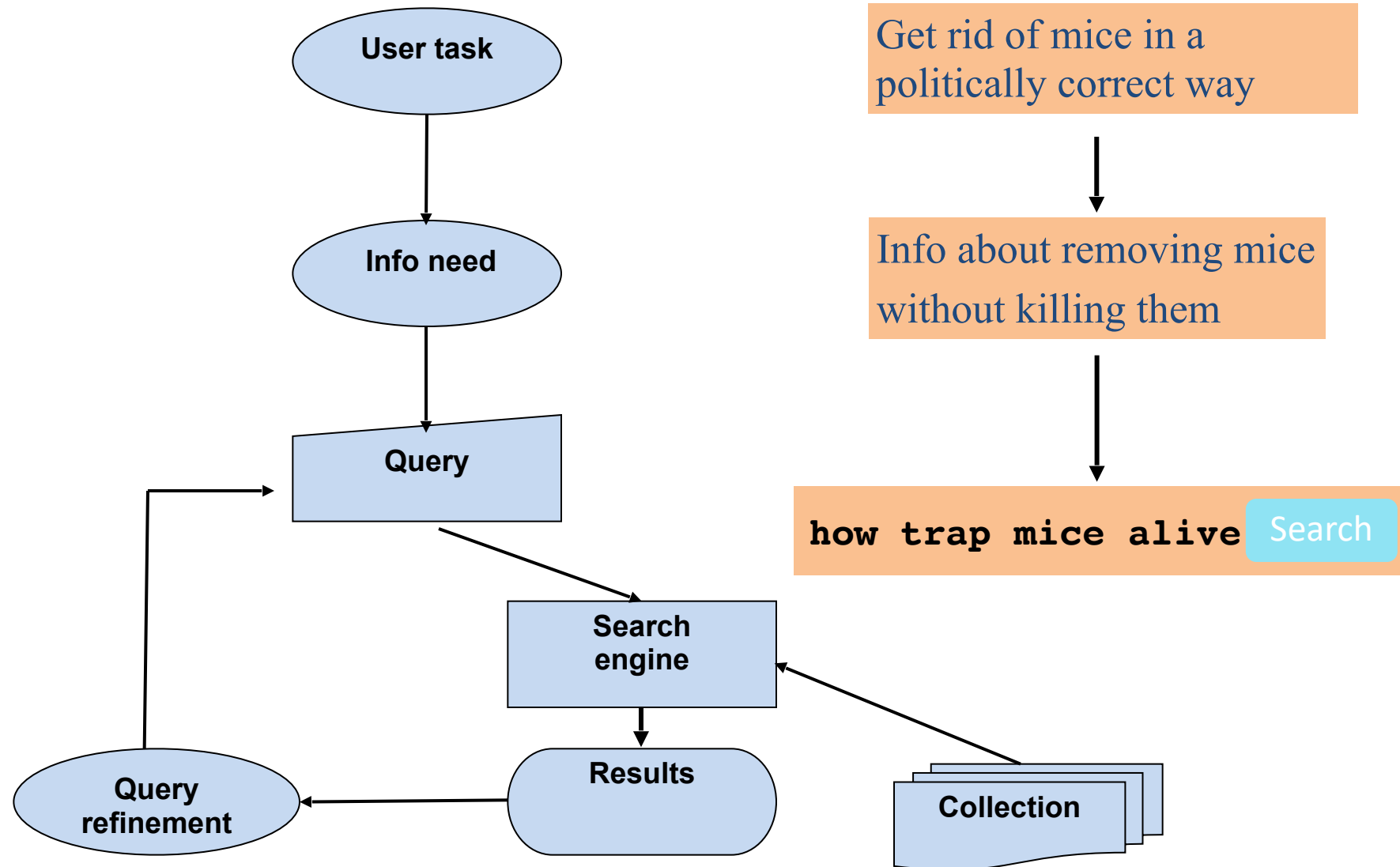
Collection

Query refinement

Get rid of mice in a politically correct way

# The classic search model
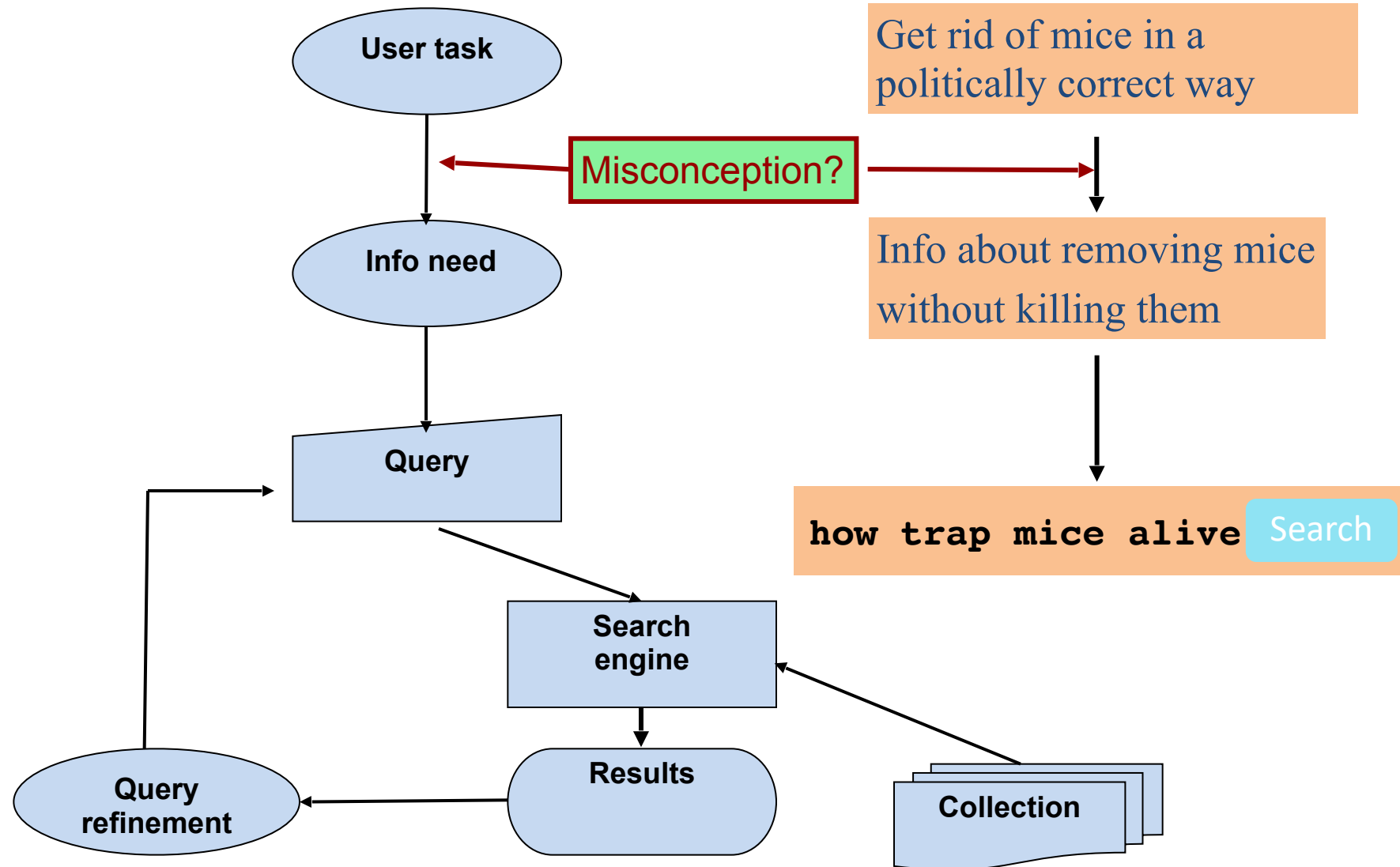
# The classic search model

# The classic search model

# The classic search model

☞       **Solution 1:  Scan each page for the string "Computer Science".**

☞    **Solution 1:  Scan each page for the string "Computer Science".**



Wait till your next life !

☞　**Solution 1:** **Scan each page for the string "Computer Science".**

Wait till your next life !

**How did Google do?**

☞    **Solution 1:  Scan each page for the string "Computer Science".**

抓狂啦！

Wait till your next life !

Google™ 谷歌

Have more than 1 billion web pages Indexed|

Google 搜索    手气不错

高级搜索
使用偏好
语言工具

10

☞    **Solution 1:  Scan each page for the string "Computer Science".**

抓狂啦!

Wait till your next life !

Have more than 1 trillion web pages Indexed

Google 搜索      手气不错

高级搜索
使用偏好
语言工具

10

☞    **Solution 2:** Term-Document Incidence Matrix

☞ **Solution 2: Term-Document Incidence Matrix**

〚**Example**〛 **Document sets**

| Doc | Text |
|---|---|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

11

☞     **Solution 2: Term-Document Incidence Matrix**

〖**Example**〗   **Document sets**

| Doc | Text |
|-----|------|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

|  | **1** | **2** | **3** | **4** |
|----|----|----|----|----|
| a | 0 | 1 | 1 | 1 |
| arrived | 0 | 0 | 1 | 1 |
| damaged | 0 | 1 | 0 | 0 |
| delivery | 0 | 0 | 1 | 0 |
| fire | 0 | 1 | 0 | 0 |
| gold | 1 | 1 | 0 | 1 |
| of | 0 | 1 | 1 | 1 |
| in | 0 | 1 | 1 | 1 |
| shipment | 0 | 1 | 0 | 1 |
| silver | 1 | 0 | 1 | 0 |
| truck | 1 | 0 | 1 | 1 |

☞     **Solution 2: Term-Document Incidence Matrix**

〖**Example**〗   **Document sets**

| Doc | Text |
|---|---|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

|  | **1** | **2** | **3** | **4** |
|---|---|---|---|---|
| a | 0 | 1 | 1 | 1 |
| arrived | 0 | 0 | 1 | 1 |
| damaged | 0 | 1 | 0 | 0 |
| delivery | 0 | 0 | 1 | 0 |
| fire | 0 | 1 | 0 | 0 |
| gold | 1 | 1 | 0 | 1 |
| of | 0 | 1 | 1 | 1 |
| in | 0 | 1 | 1 | 1 |
| shipment | 0 | 1 | 0 | 1 |
| silver | 1 | 0 | 1 | 0 |
| truck | 1 | 0 | 1 | 1 |

**silver & truck**

11

☞     **Solution 2:** **Term-Document Incidence Matrix**

〖**Example**〗   **Document sets**

| Doc | Text |
|-----|------|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

|          | **1** | **2** | **3** | **4** |
|----------|-------|-------|-------|-------|
| a        | 0 | 1 | 1 | 1 |
| arrived  | 0 | 0 | 1 | 1 |
| damaged  | 0 | 1 | 0 | 0 |
| delivery | 0 | 0 | 1 | 0 |
| fire     | 0 | 1 | 0 | 0 |
| gold     | 1 | 1 | 0 | 1 |
| of       | 0 | 1 | 1 | 1 |
| in       | 0 | 1 | 1 | 1 |
| shipment | 0 | 1 | 0 | 1 |
| silver   | 1 | 0 | 1 | 0 |
| truck    | 1 | 0 | 1 | 1 |

**silver & truck = 1010 & 1011 = 1010**

☞    **Solution 2:** **Term-Document Incidence Matrix**

〖**Example**〗    **Document sets**

| Doc | Text |
|---|---|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

|  | **1** | **2** | **3** | **4** |
|---|---|---|---|---|
| a | 0 | 1 | 1 | 1 |
| arrived | 0 | 0 | 1 | 1 |
| damaged | 0 | 1 | 0 | 0 |
| delivery | 0 | 0 | 1 | 0 |
| fire | 0 | 1 | 0 | 0 |
| gold | 1 | 1 | 0 | 1 |
| of | 0 | 1 | 1 | 1 |
| in | 0 | 1 | 1 | 1 |
| shipment | 0 | 1 | 0 | 1 |
| silver | 1 | 0 | 1 | 0 |
| truck | 1 | 0 | 1 | 1 |

**silver & truck = 1010 & 1011 = 1010**

11

☞ **Solution 3: Compact Version - Inverted File Index**

☞    **Solution 3:  Compact Version - Inverted File Index**

【Definition】  **Index** is a mechanism for locating a given term in a text.

☞ **Solution 3:** **Compact Version - Inverted File Index**

〖Definition〗 **Index is a mechanism for locating a given term in a text.**

〖Definition〗 **Inverted file contains a list of pointers (e.g. the number of a page) to all occurrences of that term in the text.**

12

☞   **Solution 3:  Compact Version - Inverted File Index**

【Definition】 **Index** is a mechanism for locating a given term in a text.

【Definition】 **Inverted file** contains a list of pointers (e.g. the number of a page) to all occurrences of that term in the text.

| Doc | Text |
|-----|------|
| 1 | Gold silver truck |
| 2 | Shipment of gold damaged in a fire |
| 3 | Delivery of silver arrived in a silver truck |
| 4 | Shipment of gold arrived in a truck |

☞ **Solution 3:** **Compact Version - Inverted File Index**

【Definition】 **Index** is a mechanism for locating a given term in a text.

【Definition】 **Inverted file** contains a list of pointers (e.g. the number of a page) to all occurrences of that term in the text.

| Doc | Text |
|-----|------|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

**Inverted File Index** →

| No. | Term | Times; Documents |
|-----|------|------------------|
| 1 | a | <3; 2,3,4> |
| 2 | arrived | <2; 3,4> |
| 3 | damaged | <1; 2> |
| 4 | delivery | <1; 3> |
| 5 | fire | <1; 2> |
| 6 | gold | <3; 1,2,4> |
| 7 | of | <3; 2,3,4> |
| 8 | in | <3; 2,3,4> |
| 9 | shipment | <2; 2,4> |
| 10 | silver | <2; 1,3> |
| 11 | truck | <3; 1,3,4> |

12

☞  **Solution 3:** **Compact Version - Inverted File Index**

【Definition】 **Index** **is a mechanism for locating a given term in a text.**

【Definition】 **Inverted file** **contains a list of pointers (e.g. the number of a page) to all occurrences of that term in the text.**

| Doc | Text |
|-----|------|
| 1 | Gold silv... |
| 2 | Shipment of g... damaged in a fire |
| 3 | Delivery of silver arrived in a silver truck |
| 4 | Shipment of gold arrived in a truck |

Index →

| No. | Term | Times; Documents |
|-----|------|------------------|
|  |  | <..; 2,3,4> |
|  |  |  |
|  |  |  |
|  |  | <..; ..> |
| 5 | fire | <1; 2> |
| 6 | gold | <3; 1,2,4> |
| 7 | of | <3; 2,3,4> |
| 8 | in | <3; 2,3,4> |
| 9 | shipment | <2; 2,4> |
| 10 | silver | <2; 1,3> |
| 11 | truck | <3; 1,3,4> |

**Inverted because it lists for a *term*, all documents that contain the term**

12

| Doc | Text |
|-----|------|
| 1 | Gold silver truck |
| 2 | Shipment of gold damaged in a fire |
| 3 | Delivery of silver arrived in a silver truck |
| 4 | Shipment of gold arrived in a truck |

| No. | Term | Times; Documents |
|-----|------|------------------|
| 1 | a | <3; 2,3,4> |
| 2 | arrived | <2; 3,4> |
| 3 | damaged | <1; 2> |
| 4 | delivery | <1; 3> |
| 5 | fire | <1; 2> |
| 6 | gold | <3; 1,2,4> |
| 7 | of | <3; 2,3,4> |
| 8 | in | <3; 2,3,4> |
| 9 | shipment | <2; 2,4> |
| 10 | silver | <2; 1,3> |
| 11 | truck | <3; 1,3,4> |

| Doc | Text |
|---|---|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

| No. | Term | Times; Documents |
|---|---|---|
| 1 | a | <3; 2,3,4> |
| 2 | arrived | <2; 3,4> |
| 3 | damaged | <1; 2> |
| 4 | delivery | <1; 3> |
| 5 | fire | <1; 2> |
| 6 | gold | <3; 1,2,4> |
| 7 | of | <3; 2,3,4> |
| 8 | in | <3; 2,3,4> |
| 9 | shipment | <2; 2,4> |
| 10 | silver | <2; 1,3> |
| 11 | truck | <3; 1,3,4> |

**How to easily print the sentences which contain the words and highlight the words?**

| Doc | Text |
|---|---|
| 1 | Gold silver truck |
| 2 | Shipment of gold damaged in a fire |
| 3 | Delivery of silver arrived in a silver truck |
| 4 | Shipment of gold arrived in a truck |

| No. | Term | Times; Documents Words |
|---|---|---|
| 1 | a | <3; (2;6),(3;6),(4;6)> |
| 2 | arrived | <2; (3;4),(4;4)> |
| 3 | damaged | <1; (2;4)> |
| 4 | delivery | <1; (3;1)> |
| 5 | fire | <1; (2;7)> |
| 6 | gold | <3; (1;1),(2;3),(4;3)> |
| 7 | of | <3; (2;2),(3;2),(4;2)> |
| 8 | in | <3; (2;5),(3;5),(4;5)> |
| 9 | shipment | <2; (2;1),(4;1)> |
| 10 | silver | <2; (1;2),(3;3,7)> |
| 11 | truck | <3; (1;3),(3;8),(4;7)> |

How to easily print the sentences which contain the words and highlight the words?

| Doc | Text |
|---|---|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

| No. | Term | Times; Documents Words |
|---|---|---|
| 1 | a | <3; (2;6),(3;6),(4;6)> |
| 2 | arrived | <2; (3;4),(4;4)> |
| 3 | damaged | <1; (2;4)> |
| 4 | delivery | <1; (3;1)> |
| 5 | fire | <1; (2;7)> |
| 6 | gold | <3; (1;1),(2;3),(4;3)> |
| 7 | of | <3; (2;2),(3;2),(4;2)> |
| 8 | in | <3; (2;5),(3;5),(4;5)> |
| 9 | shipment | <2; (2;1),(4;1)> |
| 10 | silver | <2; (1;2),(3;3,7)> |
| 11 | truck | <3; (1;3),(3;8),(4;7)> |

**Term Dictionary**

**How to easily print the sentences which contain the words and highlight the words?**

13

| Doc | Text |
|-----|------|
| 1 | **Gold silver truck** |
| 2 | **Shipment of gold damaged in a fire** |
| 3 | **Delivery of silver arrived in a silver truck** |
| 4 | **Shipment of gold arrived in a truck** |

| No. | Term | Times; Documents Words |
|-----|------|------------------------|
| 1 | a | <3; (2;6),(3;6),(4;6)> |
| 2 | arrived | <2; (3;4),(4;4)> |
| 3 | damaged | <1; (2;4)> |
| 4 | delivery | <1; (3;1)> |
| 5 | fire | <1; (2;7)> |
| 6 | gold | <3; (1;1),(2;3),(4;3)> |
| 7 | of | <3; (2;2),(3;2),(4;2)> |
| 8 | in | <3; (2;5),(3;5),(4;5)> |
| 9 | shipment | <2; (2;1),(4;1)> |
| 10 | silver | <2; (1;2),(3;3,7)> |
| 11 | truck | <3; (1;3),(3;8),(4;7)> |

**Term Dictionary**     **Posting List**

**How to easily print the sentences which contain the words and highlight the words?**

| Doc | Text |
|-----|------|
| 1 | Gold silver truck |
| 2 | Shipment of gold damaged in a fire |
| 3 | Delivery of silver arrived in a silver truck |
| 4 | Shipment of gold arrived in a truck |

| No. | Term | Times; Documents Words |
|-----|------|------------------------|
| 1 | a | <3; (2;6),(3;6),(4;6)> |
| 2 | arrived | <2; (3;4),(4;4)> |
| 3 | damaged | <1; (2;4)> |
| 4 | delivery | <1; (3;1)> |
| 5 | fire | <1; (2;7)> |
| 6 | gold | <3; (1;1),(2;3),(4;3)> |
| 7 | of | <3; (2;2),(3;2),(4;2)> |
| 8 | in | <3; (2;5),(3;5),(4;5)> |
| 9 | shipment | <2; (2;1),(4;1)> |
| 10 | silver | <2; (1;2),(3;3,7)> |
| 11 | truck | <3; (1;3),(3;8),(4;7)> |

Term Dictionary          Posting List

How to easily print the sentences which contain the words and highlight the words?

Why do we keep "times" (frequency)?

13

**Index Generator**

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Token Analyzer
Stop Filter**

14

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Token Analyzer**
**Stop Filter**

14

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Token Analyzer**
**Stop Filter**

**Vocabulary**
**Scanner**

14

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Token Analyzer
Stop Filter**

**Vocabulary
Scanner**

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

**Token Analyzer
Stop Filter**

**Vocabulary
Scanner**

**Vocabulary
Insertor**

14

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

| Token Analyzer Stop Filter | Vocabulary Scanner | Vocabulary Insertor |
|---|---|---|

**Index Generator**

```
while ( read a document D ) {
    while ( read a term T in D ) {
        if ( Find( Dictionary, T ) == false )
            Insert( Dictionary, T );
        Get T's posting list;
        Insert a node to T's posting list;
    }
}
Write the inverted index to disk;
```

| Token Analyzer Stop Filter |

| Vocabulary Scanner |

| Vocabulary Insertor |

| Memory management |

**While reading a term ……**

**While reading a term ……**

✄    *Word Stemming*

**While reading a term ……**

✂ *Word Stemming*

**Process a word so that only its stem or root form is left.**

**While reading a term ……**

✂   *Word Stemming*

 **Process a word so that only its stem or root form is left.**

〖**Example**〗          **Process**
                       **processing**      ⎫
                       **processes**       ⎬  **process**
                       **processed**       ⎭

          **says**
          **said**    ⎫  **say**
          **saying**  ⎭

**While reading a term ……**

✂ *Word Stemming*

**Process a word so that only its stem or root form is left.**

〚**Example**〛　　**Process**
**processing**　　　　　**process**
**processes**
**processed**

**says**
**said**　　**say**
**saying**

✂ *Stop Words*

15

**While reading a term ……**

✂  *Word Stemming*

 **Process a word so that only its stem or root form is left.**

 〖**Example**〗
Process
processing
processes
processed
} process

says
said
saying
} say

✂   *Stop Words*

**Some words are so common that almost every document contains them, such as "a" "the" "it".  It is useless to index them.  They are called *stop words*.  We can eliminate them from the original documents.**

**While accessing a term ……**

**While accessing a term ……**

☞ **Solution 1: Search trees ( *B- trees, B+ trees, Tries, ...* )**

**While accessing a term ……**

☞        **Solution 1:  Search trees ( *B- trees, B+ trees, Tries, ...* )**

☞        **Solution 2:  Hashing**

**While accessing a term ……**

☞     **Solution 1:  Search trees ( *B- trees, B+ trees, Tries, ...* )**

☞     **Solution 2:  Hashing**

> **Discussion 3:**
> **What are the pros and cons of using hashing, comparing to using search trees?**

**While accessing a term ……**

☞ **Solution 1:** **Search trees ( *B- trees, B+ trees, Tries, ...* )**

☞ **Solution 2:** **Hashing**

> **Discussion 3:**
> **What are the pros and cons of using hashing, comparing to using search trees?**

☝ **faster for one word**

☟ **scanning in sequential order is not possible (e.g. range searches are expensive)**

**While not having enough memory ……**

**While not having enough memory ……**

```
while ( read a document D ) {
   while ( read a term T in D ) {



      if ( Find( Dictionary, T ) == false )
        Insert( Dictionary, T );
      Get T's posting list;
      Insert a node to T's posting list;
   }
}
for ( i=0; i<BlockCnt; i++ )


```

**While not having enough memory ……**

```
BlockCnt = 0;
while ( read a document D ) {
   while ( read a term T in D ) {
      if ( out of memory ) {
         Write BlockIndex[BlockCnt] to disk;
         BlockCnt ++;
         FreeMemory;
      }
      if ( Find( Dictionary, T ) == false )
         Insert( Dictionary, T );
      Get T's posting list;
      Insert a node to T's posting list;
   }
}
for ( i=0; i<BlockCnt; i++ )
```

**While not having enough memory ……**

```
BlockCnt = 0;
while ( read a document D ) {
  while ( read a term T in D ) {
    if ( out of memory ) {
      Write BlockIndex[BlockCnt] to disk;
      BlockCnt ++;
      FreeMemory;
    }
    if ( Find( Dictionary, T ) == false )
      Insert( Dictionary, T );
    Get T's posting list;
    Insert a node to T's posting list;
  }
}
for ( i=0; i<BlockCnt; i++ )
  Merge( InvertedIndex, BlockIndex[i] );
```

17

**While not having enough memory ……**

```
BlockCnt = 0;
while ( read a document D ) {
  while ( read a term T in D ) {
    if ( out of memory ) {
      Write BlockIndex[BlockCnt] to disk;
      BlockCnt ++;
      FreeMemory;
    }
    if ( Find( Dictionary, T ) == false )
      Insert( Dictionary, T );
    Get T's posting list;
    Insert a node to T's posting list;
  }
}
for ( i=0; i<BlockCnt; i++ )
  Merge( InvertedIndex, BlockIndex[i] );
```

**Sorted**

17

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**—— Each node contains index of a subset of collection**

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**——— Each node contains index of a <u>subset</u> of collection**

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**—— Each node contains index of a <u>subset</u> of collection**

☞    **Solution 1: Term-partitioned index**

**Distributed indexing (for web-scale indexing — don't try this at home!)**

—— **Each node contains index of a <u>subset</u> of collection**

☞ **Solution 1: Term-partitioned index**



A~C          D~F          …………          X~Z

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**—— Each node contains index of a <u>subset</u> of collection**

☞ **Solution 1:** **Term-partitioned index**



**A~C**     **D~F** ············     **X~Z**

☞ **Solution 2:** **Document-partitioned index**

**Distributed indexing (for web-scale indexing — don't try this at home!)**

**—— Each node contains index of a <u>subset</u> of collection**

☞　　**Solution 1:** **Term-partitioned index**



**A~C**　　　**D~F**　　　　　　　**X~Z**

☞　　**Solution 2:** **Document-partitioned index**



**1~**　　　**10001~**　　　　　**90001~**
**10000**　　**20000**　　　　　**100000**

18

**Dynamic indexing**

**Dynamic indexing**

☞    **Docs come in over time**

        **- postings updates for terms already in dictionary**
        **- new terms added to dictionary**

☞ **Docs get deleted**

**Dynamic indexing**

☞    **Docs come in over time**

- **postings updates for terms already in dictionary**
- **new terms added to dictionary**

☞ **Docs get deleted**

**Main Index**

19

**Dynamic indexing**

☞ **Docs come in over time**

   **- postings updates for terms already in dictionary**
   **- new terms added to dictionary**

☞ **Docs get deleted**

**Main Index**

**auxiliary index**

**Dynamic indexing**

☞ **Docs come in over time**

- **postings updates for terms already in dictionary**
- **new terms added to dictionary**

☞ **Docs get deleted**

📄 **New Docs**

**auxiliary index**

**Main Index**

**Dynamic indexing**

☞ **Docs come in over time**

    **- postings updates for terms already in dictionary**
    **- new terms added to dictionary**

☞ **Docs get deleted**

📑 **New Docs**

**auxiliary index**

**Main Index**

**Search Results**

**Dynamic indexing**

☞ **Docs come in over time**

    **- postings updates for terms already in dictionary**
    **- new terms added to dictionary**

☞ **Docs get deleted**

**New Docs**

**auxiliary index**

**Main Index**

**Search Results**

**Dynamic indexing**

☞ **Docs come in over time**

     **- postings updates for terms already in dictionary**
     **- new terms added to dictionary**

☞ **Docs get deleted**

**New Docs**

**auxiliary index**

**Main Index**

**Search Results**

**When to re-index?**
**How to delete a doc?**

**Very short answer (for exam-style):**

1. **When to re-index?**
   - When the auxiliary index grows beyond a size/segment threshold or at scheduled low-traffic times, often also when there are many deletions. Then merge auxiliary and main index (and possibly rebuild).

2. **How to delete a document?**
   - Use **lazy deletion**: mark the document as deleted in a separate deletion list / bit-vector so queries ignore it, and **physically remove** its postings only when you rebuild/merge the index.

**Compression**

**Compression**

| Term |
| --- |
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

## Compression

| Term |
|------|
| ~~a~~ |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| ~~of~~ |
| ~~in~~ |
| shipment |
| silver |
| truck |

## Compression

| Term |
|------|
| ~~a~~ |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| ~~of~~ |
| ~~in~~ |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

## Compression

| Term |
|------|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ | arrivedamagedeliverfiregoldshipmentsilvertruck |

## Compression

| Term |
|---|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

**Compression**

| Term |
|------|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

21

**Compression**

| Term |
|---|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

2, 13, 32, … …, 11, …

21

**Compression**

| Term |
|------|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

2, 13, 32, … …, 11, …

21

## Compression

| Term |
|------|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

2, 13, 32, … …, 11, …

**Compression**

| Term |
|---|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡️ **arrivedamagedeliverfiregoldshipmentsilvertruck**

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

2, 13, 32, … …, 11, …

**Compression**

| Term |
|------|
| a |
| arrive |
| damage |
| deliver |
| fire |
| gold |
| of |
| in |
| shipment |
| silver |
| truck |

➡ arrivedamagedeliverfiregoldshipmentsilvertruck

**Posting List:**

computer ⟶ 2, 15, 47, …, 58879, 58890, …

2, 13, 32, … …, 11, …

☞ **Most gaps can be encoded with far fewer than 20 bits**

21

**Thresholding**

**Thresholding** 阈值 != 阀值

**Thresholding**

☞ **Document:** only retrieve the top *x* documents where the documents
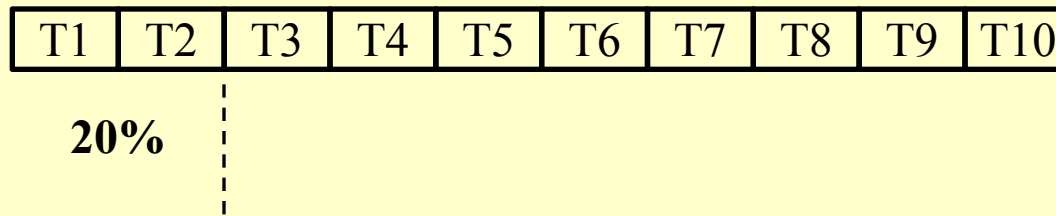
are ranked by weight

**Thresholding**

☞ **Document:** only retrieve the top *x* documents where the documents are ranked by weight

☞ Not feasible for Boolean queries

☞ Can miss some relevant documents due to truncation

**Thresholding**

☞ **Document:** **only retrieve the top *x* documents where the documents**

**are ranked by weight**

👎     **Not feasible for Boolean queries**

👎  **Can miss some relevant documents due to
truncation**

☞ **Query:** **Sort the query terms by their frequency in ascending**

**order; search according to only some percentage of the original query
terms**

**Thresholding**

☞ **Document:** only retrieve the top *x* documents where the documents

are ranked by weight

👎 Not feasible for Boolean queries

👎 Can miss some relevant documents due to
truncation

☞ **Query:** Sort the query terms by their frequency in ascending

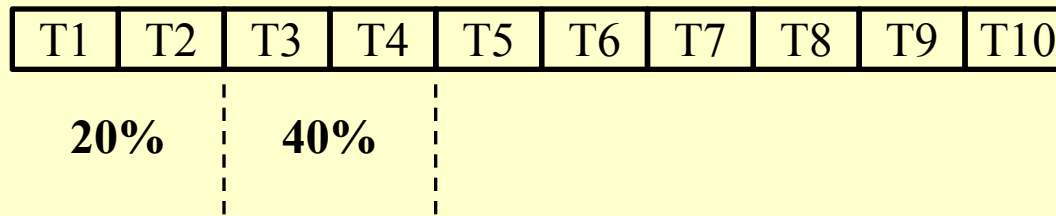order; search according to only some percentage of the original query
terms

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 |
|----|----|----|----|----|----|----|----|----|-----|

22

**Thresholding**

☞ **Document:** **only retrieve the top** *x* **documents where the documents are ranked by weight**

👎 **Not feasible for Boolean queries**

👎 **Can miss some relevant documents due to truncation**

☞ **Query:** **Sort the query terms by their frequency in ascending order; search according to only some percentage of the original query terms**

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 |
|----|----|----|----|----|----|----|----|----|-----|

**20%**

**Thresholding**

☞ **Document:** only retrieve the top *x* documents where the documents

are ranked by weight

👎 Not feasible for Boolean queries

👎 Can miss some relevant documents due to
truncation

☞ **Query:** Sort the query terms by their frequency in ascending

order; search according to only some percentage of the original query
terms

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 |
|----|----|----|----|----|----|----|----|----|-----|

  **20%**      **40%**

**Thresholding**

☞ **Document:** only retrieve the top *x* documents where the documents are ranked by weight

👎 **Not feasible for Boolean queries**

👎 **Can miss some relevant documents due to truncation**

☞ **Query:** Sort the query terms by their frequency in ascending order; search according to only some percentage of the original query terms
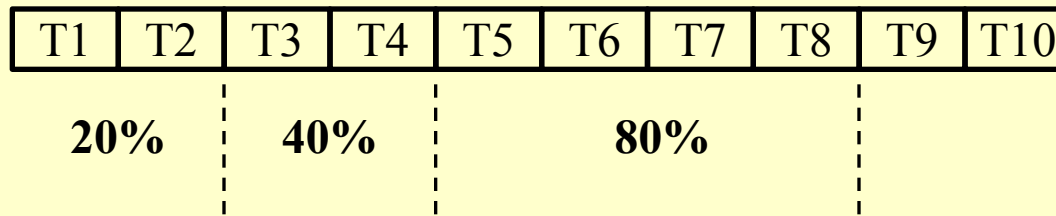
| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 |
|----|----|----|----|----|----|----|----|----|-----|

**20%**     **40%**     **80%**

**Measures for a search engine**

**Measures for a search engine**

☞ **How fast does it index**

**- Number of documents/hour**

☞ **How fast does it search**

**- Latency as a function of index size**

☞ **Expressiveness of query language**

**- Ability to express complex information needs**
**- Speed on complex queries**

**Measures for a search engine**

☞    **How fast does it index**

- **Number of documents/hour**

☞ **How fast does it search**

- **Latency as a function of index size**

☞ **Expressiveness of query language**

- **Ability to express complex information needs**
- **Speed on complex queries**

☞    **User happiness**

**Measures for a search engine**

☞ **How fast does it index**

- **Number of documents/hour**

☞ **How fast does it search**

- **Latency as a function of index size**

☞ **Expressiveness of query language**

- **Ability to express complex information needs**
- **Speed on complex queries**

☞ **User happiness ?**

**Measures for a search engine**

☞ **How fast does it index**

- **Number of documents/hour**

☞ **How fast does it search**

- **Latency as a function of index size**

☞ **Expressiveness of query language**

- **Ability to express complex information needs**
- **Speed on complex queries**

☞ **User happiness ?**

- **Data Retrieval Performance Evaluation (after establishing correctness)**
  > **Response time**
  > **Index space**
- **Information Retrieval Performance Evaluation**
  > **+ How *relevant* is the answer set?**

*Relevance* measurement requires 3 elements:

1.   A benchmark document collection
2.   A benchmark suite of queries
3.   A binary assessment of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

*Relevance* measurement requires 3 elements:

1. A benchmark **document** collection
2. A benchmark suite of **queries**
3. A binary **assessment** of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

|  | Relevant | Irrelevant |
|---|---|---|
| Retrieved | $R_R$ | $I_R$ |
| Not Retrieved | $R_N$ | $I_N$ |

*Relevance* measurement requires 3 elements:

1. A benchmark **document** collection
2. A benchmark suite of **queries**
3. A binary **assessment** of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

|  | Relevant | Irrelevant |
|---|---|---|
| **Retrieved** | $R_R$ | $I_R$ |
| **Not Retrieved** | $R_N$ | $I_N$ |

**Precision** $P = R_R / (R_R + I_R)$
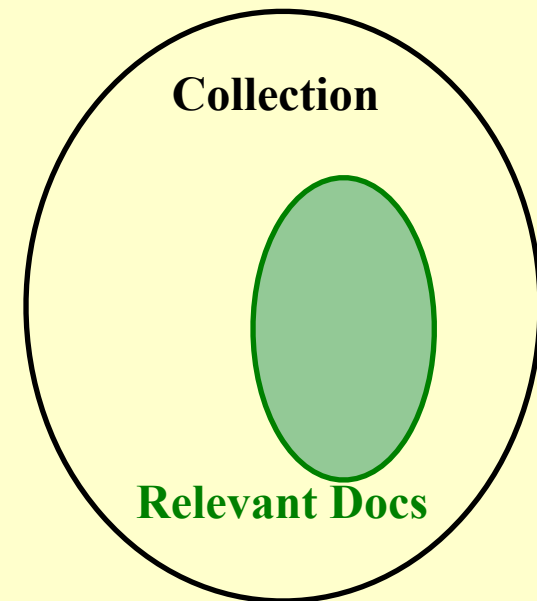
**Recall** $R = R_R / (R_R + R_N)$

*Relevance* measurement requires 3 elements:

1. A benchmark **document** collection
2. A benchmark suite of **queries**
3. A binary **assessment** of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

|  | Relevant | Irrelevant |
|---|---|---|
| **Retrieved** | $R_R$ | $I_R$ |
| **Not Retrieved** | $R_N$ | $I_N$ |

**Collection**

**Precision** $P = R_R / (R_R + I_R)$

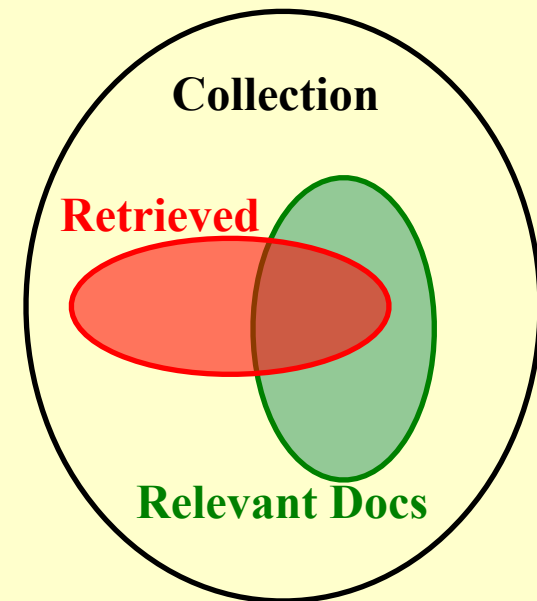**Recall** $R = R_R / (R_R + R_N)$

24

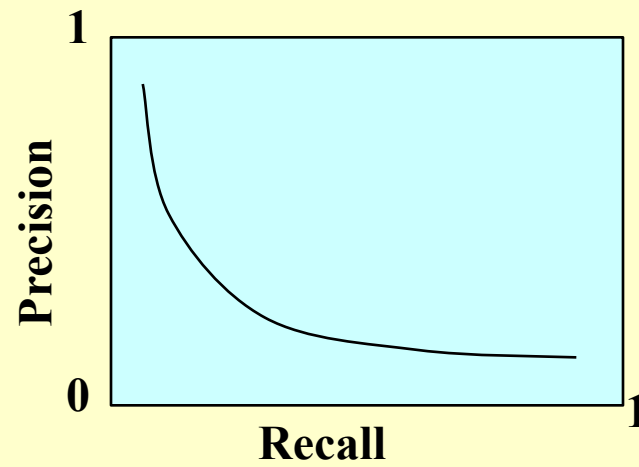*Relevance* measurement requires 3 elements:

1.  A benchmark **document** collection
2.  A benchmark suite of **queries**
3.  A binary **assessment** of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

|               | Relevant | Irrelevant |
| ------------- | -------- | ---------- |
| Retrieved     | $R_R$    | $I_R$      |
| Not Retrieved | $R_N$    | $I_N$      |

**Precision** $P = R_R / (R_R + I_R)$

**Recall** $R = R_R / (R_R + R_N)$

Collection

Relevant Docs

24

*Relevance* measurement requires 3 elements:

1. A benchmark document collection
2. A benchmark suite of queries
3. A binary assessment of either <u>Relevant</u> or <u>Irrelevant</u> for each query-doc pair

|  | Relevant | Irrelevant |
|---|---|---|
| **Retrieved** | $R_R$ | $I_R$ |
| **Not Retrieved** | $R_N$ | $I_N$ |

Precision $P = R_R / (R_R + I_R)$

Recall $R = R_R / (R_R + R_N)$

**Collection**

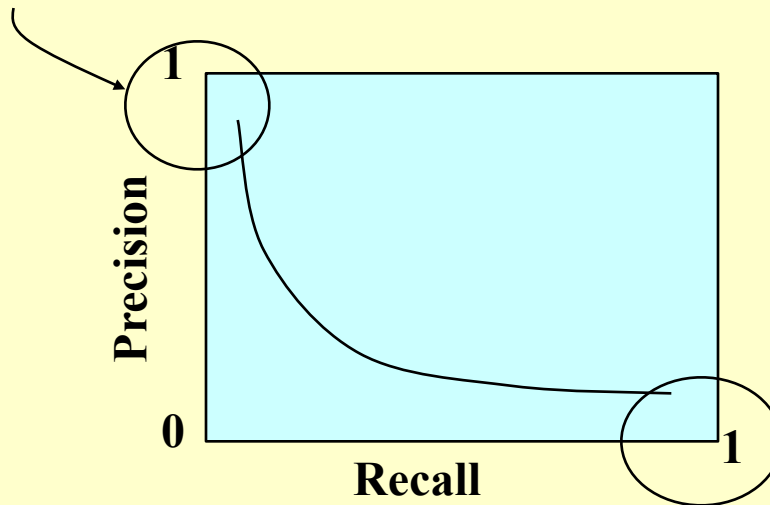**Retrieved**

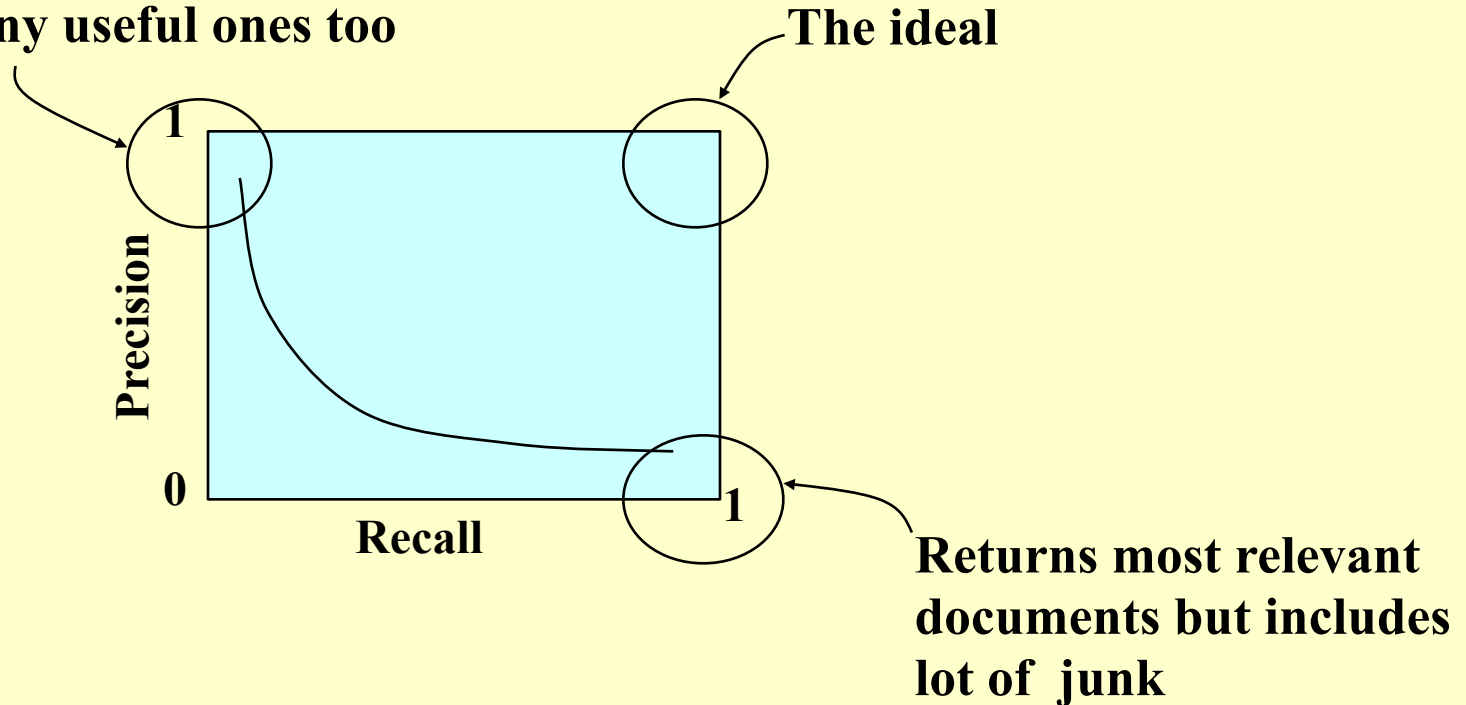**Relevant Docs**

**Returns relevant documents but
misses many useful ones too**

**Returns relevant documents but misses many useful ones too**

**Returns most relevant documents but includes lot of junk**

Precision

Recall

1

0

1

**Returns relevant documents but misses many useful ones too**

**The ideal**



Precision / Recall graph: axes labeled Precision (vertical, 0 to 1) and Recall (horizontal, 0 to 1).

**Returns most relevant documents but includes lot of junk**

25

**Discussion 4:**
**How to improve the *relevancy* of search results?**

**Discussion 4:**
**How to improve the *relevancy* of search results?**

☞ **Page Rank**

☞ **Semantic Web**

☞ **……**

# Inverted File Index

- Inverted File Index

- Take-home messages

# Take-Home Messages

- Inverted file index:

    - Mapping from items to posting lists of documents.

    - Sorted, frequency of appearance, positions.

    - Performance measures: precision, recall…

# Thanks for your attention!
# Discussions?

# Reference

Stanford CS276: https://web.stanford.edu/class/cs276/

Schütze, Hinrich, Christopher D. Manning, and Prabhakar Raghavan. Introduction to information retrieval. Cambridge University Press, 2008. Chap. 1, 4, 5, 6, 8

Zobel, Justin, and Alistair Moffat. "Inverted files for text search engines." ACM computing surveys, 2006.